

コマンドラインを使い、PDFをテキストにする

朝日新聞国際報道部・西尾 能人
yoshito.nishio+JNPC@gmail.com

0. はじめに

ふだんはマウスを使って直感的に操作しているパソコン。でも、SF映画で、「黒い画面」に命令を打ち込んで操作するシーンを見たことがあるはず。近未来的とか、プロっぽい、とかいったイメージが託されている.....のかも。

昔ならともかく、実は今でもこの「黒い画面」を使う人がいる。どんなメリットがあるのか、少しだけ体験する。その例として、PDFからテキストを抽出してみる。

下記の事前準備が難しかった方は、きょうは見るだけにして、「こういうものがある」ということだけ記憶に留めておく。大量のファイルに同じ処理をしたい、というニーズに直面したときに思い出そう。技術に詳しい人に声を掛けて、手伝ってもらえばよい。

1. このコマの目標

- Linuxのコマンドラインを知る
- コマンド操作で、大量のPDFをテキストに変換する

2. このコマでやること

- 一つのことをやる「コマンド」を組み合わせることで、複雑な処理をする
- 大量のファイルに同じ処理を適用する
- 文字情報を持っているPDFからのテキスト抽出。Popplerというツール群に含まれている「pdftotext」を使う

3. 可能なら事前準備を

ただし、**決して無理しないこと!**

自信がない場合は、詳しい方に手伝ってもらうか、インストールしないままで見学を。

3.1. Windows PCへのWSL2とUbuntuのインストール

UbuntuというのはLinuxの一種。以下の手順で、デフォルトでインストールされる。

1. **Windows**キー + **X** キーを押し、メニューから「Windows PowerShell(管理者)」を選ぶ。または、管理者権限ありで、PowerShellかcmd.exeを実行する。
2. 開いた黒い窓で、以下のコマンドを実行する。

```
ws1 --install
```

あとは、画面の指示に従う。

[参考：Microsoftのサイトの手順案内](#)

- ペンギンマークのアイコンで、Linuxの端末が開く。初回のみ、ユーザー設定が必要。

参考：Microsoftのサイトの手順案内

- もし日本語環境になっていないようなら、そのときのみ、以下を実行

★ここが今でも必要か不明。どなたか教えて下さい

```
sudo apt install language-pack-ja
sudo update-locale LANG=ja_JP.UTF-8
```

3.2. Windows PCへの「pdftotext」のインストール

- Linuxのターミナルで、以下のコマンドを実行。

```
sudo apt install poppler-utils
```

3.3. Macへの「pdftotext」のインストール

- Macのターミナルで、以下のコマンドを実行（Homebrewを使っている場合）。

```
brew install poppler
```

4. このコマでは扱わないこと

- 文字情報を持たない「画像」状態のPDFの処理

こちらは、スポーツ報知・田中さんら、本日、ほかの講師からご紹介があります。

- 文字情報を持っていて、かつ表形式の場合は「[tabula](#)」がお勧め。ブラウザを使うが、スタンドアローンで動作している
- Windows PCにLinux環境を構築する方法
- pdftotext(Poppler)のインストールがうまく行かない場合の対処

いずれも、社によってセキュリティ周りの設定やOSのバージョンが異なり、いっせいにご案内するのは難しいため。

「コマンドラインを使えばできるかも」ということが頭の片隅に残れば、必要になったときに、社内の詳しい方に相談できる。とりあえずは、それで十分。

5. フォルダの準備

練習用のファイル一式を、GitHubに置いてある。<https://github.com/nishioWU/JNPC2205>からダウンロードを。

右上の緑の「code」から「Download ZIP」を選ぶ。ダウンロード後は必ず解凍し、デスクトップなど適当な

場所に「JNPC2105-main」というファイルごと移しておく。
このコマでは、このフォルダの中だけで作業する。

6. コマンドラインを試す

6.1. Linux 端末の起動

Windows PCでは、デスクトップに準備した「JNPC2105-main」のフォルダを、**Shift**キーを押しながら右クリックする。すると、普段の右クリックでは出てこない「Linux端末をここに開く」という項目が見えるので、それを選ぶ（ここで開く、に直っているかも）。「bash」というものが立ち上がる。

Mac の場合は、ターミナルを起動する。フォルダを右クリックし、「サービス」から「ターミナル」を選んでもよい。「zsh」というものが立ち上がる。

bash と zsh で少し挙動が違う部分があるが、本日扱う範囲では、大差ない。

6.2. 使い方の基本

黒い画面に文字を打ち込み、最後に**Enter**キーを押すと、コマンドが実行される。きちんと実行された場合には、特段のメッセージは表示されない。また、危険な操作の場合も、確認・念押しはされずに、実行されてしまうので、要注意。

コマンドの履歴は**↑**キーでたどれるので（行き過ぎて戻るときは**↓**）、少し打ち替えて実行する場合に便利。打ち替える必要がなく、同じコマンドをもう一度実行する場合は**!!**。

ファイル名やコマンドは、ほかと区別のつく部分まで入力して**Tab**キーを押せば、自動補完される。区別がつかない場合には、音が鳴るので、もう少し先まで打ち足せばよい。または**Tab**を2度押すと、候補の一覧が表示される。

なお、Linuxの世界では、Windowsと違って大文字と小文字は区別され、別物として扱われるので注意。ファイルが見つからない旨のエラーが出るときは、たいてい、これが原因。

入力を中断したいときは「**Ctrl** + **C**」。実行中のコマンドをとめたいときも、同じく「**Ctrl** + **C**」。画面をいったんクリアしたいときは、「**Ctrl** + **L**」。

6.3. コマンドとオプション

まず、今日使うフォルダの構成を見てみる。「tree」コマンドを使う。黒い画面に

```
tree
```

と入力して、**Enter**を押すと実行される。見慣れたGUIではないが、ファイルとフォルダが区別して表示されているのが分かるだろうか。

あるフォルダ内のファイル一覧を見るには、「ls」コマンドを使う。やはり、黒い画面に

```
ls
```

と入力し、**Enter**を押すと実行される。

今いる場所から別のフォルダに移動するときは、

```
cd フォルダ名
```

逆に元のフォルダに戻るときには、

```
cd -
```

と入力する。1 階層上に行くときには、

```
cd ..
```

とする。よく使う。「.」1 個は今いるフォルダという意味で、2 個だと 1 階層上を指す。

動作を少し変えたいときには、ハイフンでオプションを足す。

例えば、ls コマンドでファイル名が縦に 1 つずつ表示されるようにしたければ、半角数字の「1」を添えて

```
ls -1
```

とする。

同じ意味のオプションでも、長短 2 通りの指示の出し方があることが多い。短いオプションのときはハイフン 1 つ、長いオプションのときはハイフン 2 つが基本だが、長いオプションでもハイフン 1 つのときも稀にある。

詳しい使い方を見たいときは、以下のように打つ。

```
ls --help
```

ファイルを 1 つずつ縦に並べるなら「-1」（数字のイチ）、詳しい情報を見たいときには「-l」（longのエル）、コンマで区切ってほしいければ「-m」、種類別にしてほしいときは「-X」（大文字のエックス）、タイムスタンプ順なら「-t」、順番をひっくり返したいときは「-r または --reverse」、とある。

```
ls- 1  
ls -r
```

などと入れてみよう。

オプションは組み合わせることもできて、

```
ls -ltr
```

とすれば、「1 つずつ」「タイムスタンプ」「逆順」なので、最新のファイルから順に表示することができる。

6.4. よく使うコマンド

コマンドは無数にあるが、まず使うのは

- exit (黒い画面を閉じて、GUIの世界に戻る)
- ls (今いるフォルダのアイテムを表示する)
- less (テキストファイルを1ページずつ表示する。ページ送りはスペースキー、次のファイルへの移動は→、終了はq。ファイル内の検索もできる。複数のファイルを同時に開いた場合は、次のファイルへの移動は:n、前のファイルに戻るときは:e)
- cat (複数のファイルをつなげる。または、テキストファイルを表示する)
- man (コマンドのマニュアルを表示する)
- cp (ファイルをコピーする)
- mv (コピーせずに移動する。または名前を変更する)
- grep (テキストファイルを検索する。「正規表現」というものが使える)
- sort (テキストをソートする)
- uniq (ソート済みのものについて、ダブりを削除する)
- wc (スペースで区切られた単語数や、行数を数える)
- cd (別のフォルダに移動する)
- mkdir (フォルダを新しく作る。エクスプローラーやファインダーで作っても同じこと)
- find (条件にあうファイルを探す)
- xargs (検索結果などを「引数」として次のコマンドに渡して実行する)

あたり。その都度、使い方を調べて入力すればよい。私もよく使うオプション以外は頭に入っていない。

6.5. パイプライン、リダイレクトとワイルドカード

個々のコマンドは、単一のことをしっかりやる、という設計思想で作られている。回転寿司の一皿みたいなもので、いざ、まとまった作業をしようすると、物足りない。そこで、必要に応じて複数のコマンドをつなげて使う。Aで処理した結果をBに渡して処理させ、Cというファイルに書き込む、というように。

あるコマンドを実行した結果を、別のコマンドに渡すときには、| (パイプ、パイプライン) を使う。Shiftキーを押しながら¥マークのキーも押すと、出る。

別のコマンドに渡すのではなく、ファイルに書き込むときには、>か>> (リダイレクト) を使う。1個なら上書き、2個だと追記になる。存在しないファイル名を指定すると、その名前で新たにファイルが作られる。

それから、ファイル名の指定にはワイルドカードが使える。「*」は任意の文字列 (文字なしも含む)、「?」は任意の1文字。このほかにもあるが、複雑なので、とりあえず2つだけ覚えておけばよい。

[参考：Linuxのワイルドカードと正規表現についてのサイト \(ためになるが難しい\)](#)

コマンドを複数つなげて使い、ワイルドカードを活用することで、大量のファイルを一括処理することができる。それが、あえてマウスを使わずにコマンドを使う理由。たくさんの文書やデータと格闘するときに、役立つことが多い。

6.6. grepで絞り込み、wcでカウント

では、練習のため、cdコマンドで「TEXT」フォルダに移動を。例えば、

```
grep -n 東区 text*.txt > result.txt
```

こうすると、①今のフォルダの中にある「text」で始まる名前のテキストファイルすべてから②「東区」という文字列を含む行を探して③「result.txt」というファイルに書き込むーという命令になる。result.txtというファイルがなければ作るが、もしあったら、上書きされてしまう。上書きでなく追記したいなら、「>」を「>>」に変えてやる。

grepの「-n」は、行番号を振るオプション。

今は検索対象のファイルはたった4つしかないが、ワイルドカードで指定しているので、いくつに増えても同じ命令で事足りる。

ところで、「text」と全部打たずに、「t」だけ打ってTabを押せば、自動補完されることをお忘れなく。そのとき警告音が鳴るのは、「text」で始まるファイルが複数あるのでこの先は特定できませんよ、と知らせてくれているわけ。

こんどは「南区」を検索して、先ほどのファイルに追記してみよう。↑で今のコマンドを呼び出して編集してやる。「東区」を「南区」に、「>」は「>>」に直す。

```
grep -n 南区 text*.txt >> result.txt
```

東と南がどちらが多いか数えることもできる。やはり、使ったコマンドを呼び出して編集して

```
grep 東区 text*.txt | wc -l
```

と、

```
grep 南区 source*.txt | wc -l
```

を実行して比べてみよう（わざわざ「-n」を削らなくてもよい）。最後につけたオプションは、lineのエル。語数ではなく行数をカウントする、という指示。

エラーしない限り特段のメッセージが出ないのは、このようにパイプでつながたり、ファイルに書き込んだりする際に、そのほうが便利だからだ。コマンドの実行結果とエラーメッセージは別物として管理されているため、エラーメッセージだけをログファイルに追記する、といったこともできる。

6.7. ファイルを移動する

では、次に進むために、もう不要になったファイルを掃除する。trashというフォルダを作って

```
mkdir trash
```

そこに、sで始まるファイル全てを移動させてしまおう。

```
mv s* trash
```

result.txtももう不要なので、

```
mv result.txt trash
```

で移動する。ファイル名は、全部打たずに自動補完するのが楽。

7. pdftotext を試す

7.1. 自分の居場所を移動する

今度は、「PDF」フォルダの中の「source」フォルダに移動を。いったん上の階層に行ってから、また降りていくようにする（横に移動することはできない）。1 階層上に行くには、

```
cd ..
```

先ほどまでいたフォルダに戻る、という指示の仕方なら

```
cd -
```

で親フォルダ戻り、そこで

```
cd PDF  
cd source
```

で降りていく。打ち間違いを防ぐためにも、適宜Tabキーを押して自動補完しよう。一度に

```
cd PDF/source
```

でもよい。

7.2. pdftotextの使い方

いよいよ、pdftotextを使う。これは、Popplerというツール群の中のコマンドで、文字情報を持っているPDFから、テキストだけを抜き出してくれる。

```
pdftotext ファイル名
```

とすると、「abc.pdf」というファイル名のPDFを、拡張子だけ変えた「abc.txt」というテキストファイルに変換する。

参考：Poppler 公式サイト

pdftotextの使い方は、以下が詳しい。

参考：西村めぐみさんの記事

西村さんの著書「Linux+コマンド入門——シェルとコマンドライン、基本の力」は、pdftotextに限らず、Linuxのコマンド全般でも役に立つ。ただし、オプションは主なものしか載っていないので、一度は「コマンド名 --help」で、使いたいコマンドのヘルプを見るべきだ。

pdftextコマンドの出力結果をテキストファイルに書き込むのではなく、画面に出すだけでよければ、

```
pdftotext ファイル名 -
```

のように、最後にハイフンを1つ付ける。ページの区切りで改行されるのが煩わしければ、「-npgbrk」オプションを付ける。画面をさーっと流れて行ってしまう長い文書なら、

```
pdftotext ファイル名 - | less
```

のように、パイプでlessに渡してやる。

「PDF_番号.pdf」というファイルがいくつかあるので、ちゃんとテキストに変換されるか、見てみよう。くどいけれど、ファイル名入力はなるべく自動補完で。

7.3. pdftotextの「文法」の注意点

ところで、一つずつではなく、全部まとめて調べてみよう、と考えて

```
pdftotext PDF_*.pdf -
```

とか、さらに汎用的にPDFを全部捕まえると考えて

```
pdfttotext *.pdf -
```

を試みた方はいないだろうか。私も同じことを考えた。

しかし、コマンドのヘルプを見ると、pdftotextの文法は「pdftotext [options] []」となっている（[]内は省略可能）。ファイル名は1つまたは2つしかコマンドの後には置けず、1つ目が変換すべきPDF、2つ目がある場合はそれは抽出したテキストの書き込み先、という仕様だ。なので、「PDF_*」の条件に当てはまるファイルがたまたま1つなら正常動作するが、偶然2つあった場合は、2つ目は変換元ではなく書き込み先になってしまう（すでにファイルが存在している場合は上書きされる。catコマンドで中を見るか、拡張子を.txtに変えてWindowsのメモ帳で読むかして、確かめてみて）。ファイルが3つ以上見つかった場合はエラーして動かないので、ヘルプ画面が表示される、ということになる。

では、ファイルを1つずつ処理するしかないのか？ そもそも、1つずつテキストに変換するのでよいなら、Adobe Acrobat Readerでできる。文書を開いた後、「ファイル」→「その他の形式で保存」→「テキスト」を選べば済む。

なのにわざわざコマンドラインを使う理由は、大量のファイルを一気に処理したいからだった。どうすればよいだろうか。

7.4. おまけ : Ubuntuだけの楽なやり方

Windows PCのWSL2でデフォルトでインストールされるUbuntuの場合、pdftotextをインストールすると、lessコマンドでいきなりPDF内のテキストが読めるようになる。気を利かせて、そういう仕様になっているようだ(時事通信・川上さん、ご教示多謝です)。lessは一度に複数のファイルを指定しても問題なく動くので、こうすればよい。

```
less 名前の目印*.pdf > 出力ファイル名.txt
```

この1行だけで、今いるフォルダの中にある、ある文字列で始まる名前のpdfファイルすべてを、テキストに変えて一本につなげ、出力先のテキストファイルに保存する、という指示になる。具体的には、

```
less *.pdf > extracted_text.txt
```

などとコマンドを打てばよい。Ubuntuを使っている方には朗報。

8. 大量ファイルの一括処理

8.1. xargsでコマンドを組み合わせる

そこで、xargsというコマンドを使うことにする。これは、直前のコマンドの結果を、次のコマンドの好きな位置に埋め込むのに使える。pdftextに、一度に1つずつファイル名を渡すようコントロールするのに、うってつけだ。「-i」というオプションに、@でも{}でも何でもよいので目印の記号を付けて、こんなふうを使う。

```
前のコマンド | xargs -i{} 次のコマンド 入れたい場所に{}
```

この例では、{}を書いた場所に、「前のコマンド」の実行結果が1つずつ入る。それが尽きるまで、「次のコマンド」を繰り返して実行する。「-i」の後の記号を省略した場合は、{}になる。だから、pdftotextに1つずつファイルを渡して、抽出したテキストを画面に表示するには、

```
ファイルを探すコマンド | xargs -i{} pdftotext {} -
```

とすればよい。

8.2. findでファイルを探す

PDF形式のファイルを探すには、findコマンドを使うことにする。lsとは違い、今いるフォルダだけでなく、それより下の階層をすべて探してくれる。

```
find -name '*.pdf'
```

のように使う（-nameオプションはハイフン1つ）。その後は、「こんな名前のファイル」という検索条件を、シングルクォート（アポストロフィ。Shift+7で出る）で囲む。もし、並び順を辞書順ではなく、自然な数字順にしたければ、sortに渡してやる。lsの場合は小文字で「-v」とするのだが、sortは大文字で「-V」とするので、こうなる。

```
find -name '*.pdf' | sort -V
```

8.3. 1行プログラム完成

これを、先ほど試したものにつなげて、

```
find -name '*.pdf' | sort -V | xargs -i{} pdftotext {} -
```

こうなる。あと一息。画面に流すのではなく、全部をまとめて1つのテキストファイルに保存しておくなら、ファイルへのリダイレクトを足す。「joined_text.txt」というファイル名にするなら、以下のようになる。ついでに、長い文書の場合にPDFのページ区切りを取り除く、-npgbrkというオプションも入れた。

```
find -name '*.pdf' | sort -V | xargs -i{} pdftotext {} -npgbrk - >
joined_text.txt
```

段階的に試しながら、コマンドを履歴から編集して打ち足してきた。いきなりこれを見せられたら、見ないふりをしたくなる呪文だが、今なら何とか分かる、のではないだろうか。たった1行だけではあるけれど、簡単なプログラムがこれで完成だ。お疲れさま！

xargsの「-i」オプションの後に記号が書いてない場合は「{ }」と見なしてくれるので、もしも一から打ちこむ場合には、抜いてしまっても構わない。するとこうなる。

```
find -name '*.pdf' | sort -V | xargs -i pdftotext {} -npgbrk - >
joined_text.txt
```

出力場所を変えたいとき、例えば「今いるフォルダと同じ階層にある、outputというフォルダ」の下にしたときには、

```
find -name '*.pdf' | sort -V | xargs -i pdftotext {} -npgbrk - >
../output/joined_text.txt
```

とする。「..」で1階層上に上がり、その下にある「output」フォルダに降りて来る、という指定の仕方をしている。

PDFなら何でも、ではなく、テキスト抽出の対象とするPDFを絞り混みたいときは、findの後の' 'の中を変えればよい。

8.4. PDFごとに変換したいとき

さて、全体をひとつのテキストファイルにまとめるのではなく、PDFごとに同名（拡張子は「.txt」に変わる）のテキストファイルに出力するのだったら、どうすればよいか。ファイルを検索するまでは一緒だが、ソートは不要になる。

もとのPDFがあるのと同じフォルダに、それぞれテキストファイルを書き出すとすると、こうなる。

```
find -name '*.pdf' | xargs -i pdftotext -nopgbrk {}
```

8.5. 画像を抜き出す

余談ながら、Popplerをインストールすると、pdftotextだけでなく、PDFから画像を抜き出すpdfimagesなどほかのコマンドも使えるようになる。「日本記者クラブ.pdf」というPDF中の画像を抜き出し、png形式でないものは変換した上で、「JNPC-xxx.png」というファイル名で保存したい場合は、

```
pdfimages -png 日本記者クラブ.pdf JNPC
```

とする。

9. 終わりに

9.1. 枯れた技術ゆえの強み

大量のファイル进行处理する場合、最初に対象のファイルを全部まとめて1つにして一度だけ処理する手もあるだろうし、ファイルごとに処理してから最後にまとめるという手順も考えられる。取り組んでいる課題に合わせて、効率がよく、手戻りしない組み立てを考えればよい。

linuxのコマンドは、先人の検証が十分になされている「枯れた技術」なので、技術に詳しい人に訊くか、ネット検索すれば、いろいろと方法はヒットするはずだ。

9.2. 質問は財産、Slackで共有を

もちろん、Slackのチャンネルに質問を投稿いただければ、私のレベルでは分からないことでも、一緒に学び、仕事に使っている仲間が助けてくれることと思う（なので、メールよりはSlackのほうが効率よいはず）。「ここが分からない」と思ったところ、「こうしたいのに」と感じたところは、多くの方に共通するはずで、質問と解決法はみんなの財産になる。遠慮は無用！

お疲れさまでした。